







Learning to Predict Navigational Patterns From Partial Observations

Robin Karlsson , *Student Member, IEEE*, Alexander Carballo , *Member, IEEE*, Francisco Lepe-Salazar , Keisuke Fujii , *Member, IEEE*, Kento Ohtani , and Kazuya Takeda , *Senior Member, IEEE*

Abstract—Human beings cooperatively navigate rule-constrained environments by adhering to mutually known navigational patterns, which may be represented as directional pathways or road lanes. Inferring these navigational patterns from incompletely observed environments is required for intelligent mobile robots operating in unmapped locations. However, algorithmically defining these navigational patterns is nontrivial. This letter presents the first self-supervised learning (SSL) method for learning to infer navigational patterns in real-world environments from partial observations only. We explain how geometric data augmentation, predictive world modeling, and an information-theoretic regularizer enable our model to predict an unbiased local directional soft lane probability (DSLPL) field in the limit of infinite data. We demonstrate how to infer global navigational patterns by fitting a maximum likelihood graph to the DSLPL field. Experiments show that our SSL model outperforms two SOTA supervised lane graph prediction models on the nuScenes dataset. We propose our SSL method as a scalable and interpretable continual learning paradigm for navigation by perception.

Index Terms—Vision-based navigation, semantic scene understanding, continual learning, learning from experience, motion and path planning.

I. INTRODUCTION

MOBILE robots perform tasks that involve traversing an environment. To navigate rule-constrained structured environments robots are required to correctly perceive and interpret the environment. This problem is called scene understanding. Navigational patterns, or directional pathways, are a core component of understanding how to traverse structured environments [1]. In particular, efficient and safe multi-agent navigation depends on each agent following mutually known navigational patterns. The patterns can be defined by explicit

Manuscript received 20 April 2023; accepted 18 June 2023. Date of publication 3 July 2023; date of current version 25 July 2023. This letter was recommended for publication by Associate Editor S. Leonard and Editor P. Vasseur upon evaluation of the reviewers' comments. This work was supported by JST SPRING under Grant JPMJSP212. (*Corresponding author: Robin Karlsson.*)

Robin Karlsson, Keisuke Fujii, and Kento Ohtani are with the Graduate School of Informatics, Nagoya University, Nagoya 464-8603, Japan (e-mail: karlsson.robin@g.sp.m.is.nagoya-u.ac.jp; fujii@i.nagoya-u.ac.jp; ohtani.kento@g.sp.m.is.nagoya-u.ac.jp).

Alexander Carballo is with the Department of Electrical, Electronic and Computer Engineering, Gifu University, Gifu 501-1193, Japan (e-mail: alex@gifu-u.ac.jp).

Francisco Lepe-Salazar is with Ludolab, Colima 28000, México (e-mail: flepe@ludolab.org).

Kazuya Takeda is with the Graduate School of Informatics, Nagoya University, Nagoya 464-8603, Japan, and also with the TIER IV, Nagoya 450-6627, Japan (e-mail: kazuya.takeda@tier4.jp).

Digital Object Identifier 10.1109/LRA.2023.3291924

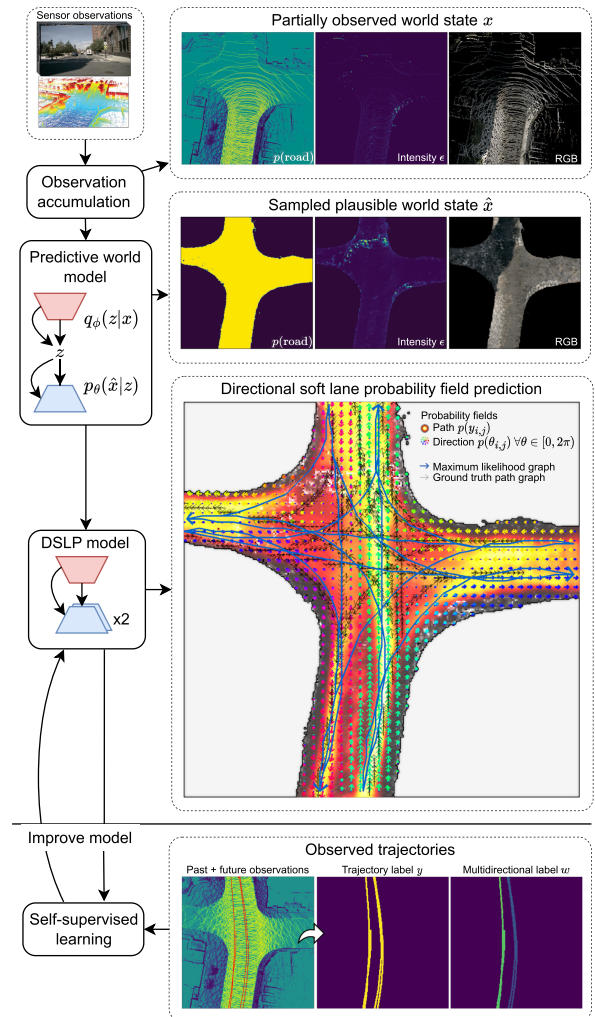


Fig. 1. Method accumulates sensor observations into a common metric vector space representing the partially observed world state x . A predictive world model samples a set of diverse plausible complete world states \hat{x} . The directional soft lane probability (DSLPL) model predicts two probability fields; the agent traversal probability $p(y_{i,j})$ and a multimodal directional probability distribution $p(\theta_{i,j})$ for each point (i,j) . A fitted maximum likelihood graph corresponds to global navigational patterns. The DSLPL model can learn navigational patterns from observed trajectories representing only a subset of all plausible trajectories.

rules or be derived from social conventions and emergent behavior. However, learning to infer navigational patterns for complex environments based on observable features is difficult due to regional variation and noise including varying or missing surface markings, geometries, and materials.

Current methods for spatial navigation can be categorized into mapping- and learning-based approaches. The mapping approach [2] avoids the problem of automatized understanding of environments by encoding human knowledge in the form of lane maps and localizing the system within these maps. Creating a priori navigation maps is a conceptually simple, interpretable, and predictable way to safely navigate environments. In practice, this approach is difficult to scale up, as map creation, maintenance, and verification are costly in terms of human labor, typically limiting application to small predetermined environments. Additionally, dynamic navigational behavior like correctly avoiding parked cars or debris cannot be a priori encoded in static maps. The learning approach involves training a model to infer navigational patterns based on environmental context. Some methods learn implicit patterns as part of accomplishing the primary task [3], [4], [5]. Other methods learn explicit patterns but require ground truth lane maps for training [6], [7]. Methods learning from observational data alone are promising scalable solutions to infer navigational patterns, as driving data can be obtained at a low cost. However, the real-world performance of existing methods is fragile and unpredictable in complex environments and lacks interpretability.

The human visual system comprises two subsystems [8], [9], [10]. The vision-for-perception system located in the ventral stream processes information in a slow, top-down manner to create perceptual representations from ambiguous or incomplete visual input by leveraging visual and semantic memory [9]. These representations support conscious mental processes such as recognition, visual thought, and planning. The vision-for-action system located in the dorsal stream processes information in a real-time, bottom-up manner to perceive the entire environment and infer behaviorally-relevant visual affordances, including cues for spatial navigation [11].

In this letter we present a self-supervised method for learning to infer navigational patterns from real-world partial observations as required for traversing unmapped real-world environments. Our approach is inspired by the biological dorsal visual pathway [9] and endows artificial intelligent agents with a functionally similar self-improving system that learns to infer visual affordances for spatial navigation [12].

The model learns general contextual environment features that explain observed trajectories, and can thus infer navigational patterns for newly encountered environments. Learning from observed trajectories means learning from only a subset of all plausible trajectories. We propose an information-theoretic regularizer to overcome the problem of false negative traversal observations resulting from partial observations. Our model combines complementary aspects of mapping- and learning-based approaches. It also produces an interpretable representation akin to maps. Lastly, this model improves with additional experience akin to continual learning [13] while avoiding catastrophic forgetting by retaining a replay buffer of past experiences [14].

We identify the navigational pattern prediction problem based on static environmental context as a sub-problem of the general dynamic agent behavior prediction problem. The main difference is that we do not consider the influence of dynamic objects such as parked cars and red traffic lights, or predict the movement of particular agents. While both problems can be solved through the same framework, we choose to remove dynamic object information from the input representation in order to objectively compare performance against ground truth lane graph methods.

While we perform experiments in a real-world urban road environment our method is applicable in any general structured environment.

The contributions of our letter are three-fold:

- A self-supervised approach for learning to predict unbiased traversability probability maps from real-world partial positive-only observations using a principled hyperparameter-free information-theoretic regularizer.
- Experimentally show that our method improves with additional observations and achieves better performance than recent state-of-the-art (SOTA) supervised methods.
- Experimentally verify that leveraging a predictive world model [15] and geometric data augmentation [16] improves real-world performance.

The rest of the letter is organized as follows. Section II reviews the SOTA and contrasts it with our work. Section III explains how partial observations are transformed into complete world states used as model input. Section IV explains our method and model implementation. Section V explains the experiment setup. Section VI present experimental results. Section VII concludes the letter by discussing limitations and future improvements to our method.

II. RELATED WORK

Path prediction: Recent works present methods to predict multimodal paths for specific actors. Salzmann et al. [17] and Baumann et al. [18] trains a convolutional neural network (CNN) on bird's-eye-view (BEV) environment representations to predict a dense map representing valid ego-vehicle paths using a weighted dense classification error and future ego-vehicle trajectories. Barnes et al. [19] trains a CNN on perspective images with self-supervised labels generated from driving data. Ort et al. [20] fuses high-level navigational guidance from a coarse map with path generation reflecting the observed environment. Casas et al. [21] optimizes a model to predict an environment map and possible paths for the ego-agent based on images and point clouds using a ground truth lane map as supervision. Prez-Higueras et al. [22] trains a CNN model to predict a multimodal path affordance map between any two points to be used as a prior for an RRT* path planner [23]. Kitani et al. [24] trains a Hidden Parameter Markov Decision Process (HiPMDP) model using inverse reinforcement learning and observation data. Ratliff et al. [25] presents an imitation learning approach that maps input features to a cost map based on example paths. Our approach expands on prior works by learning to predict all plausible navigational patterns in the environment independently of observed agents without depending on ground truth maps for supervision.

Lane graph and map prediction: Homayounfar et al. [26] trains a recurrent neural network (RNN) model to predict polylines as road lanes in highway road scenes using ground truth lane maps. An extension [27] introduces forking and merging lane topologies. Guo et al. [28] predicts 3D road lanes from perspective images using ground truth annotations. Zürn et al. [6] trains a Graph-RCNN model to predict lane anchors and edges using images and point clouds with ground truth lane map supervision. Can et al. [7] trains a transformer model to detect lane segments from images and subsequently connected into lane graphs. Zhang et al. [29] trains a three-stage network using ground truth map supervision to predict a dense lane map and subsequently predict keypoints used to generate the graph.

Mi et al. [30] presents a hierarchical coarse-to-fine approach to train an attention graph neural network to generate road lane graphs. Karlsson et al. [16] presents a self-supervised method to train a directional soft lane affordance (DSLAs) map from single trajectories. A follow-up work [31] shows how to generate discrete road lane graphs by searching for connected paths in the DSLAs map using the A^* algorithm. Our method is a scalable approach to predict lane graphs from partial observations without requiring ground truth lane map annotations and yet achieve better performance than supervised baselines [6], [7]. This work extends [16], [31] by introducing a principled regularizer, a sampling-based maximum likelihood graph generation method, and demonstrates the approach on real-world data.

Another line of works consider the problem of predicting a structured semantic representation of the environment akin to human-annotated HD maps [2] from sensor observations and ground truth maps. Li et al. [32] trains a multimodal network to predict dense maps subsequently post-processed into vectorized representations of map elements. An extension [33] directly predicts vectorized map elements. Liao et al. [34] presents a transformer model trained end-to-end to predict vectorized map elements from camera images. Shin et al. [35] presents an attention graph neural network approach. Ort [36] presents a model-based approach to fit parametric map elements according to observations and prior map information. Our approach is complementary as it provides explicit navigational patterns based on an environment representation.

End-to-end learning for autonomous vehicles: Originally proposed by Pomerleau [37] and more recently repopularized by Bojarski et al. [3], the end-to-end learning paradigm aims to learn a driving model or policy mapping perception to control by optimizing for an extrinsic goodness objective. Imitation learning approaches [3], [4], [5] learn a policy that results in similar behavior as expert examples. Reinforcement learning (RL) approaches [38] optimize a policy to maximize an extrinsically defined reward such as time-to-human-override. Recently, approaches learning an explicit predictive world model [39], [40] show that robust policies can be learned from expert observation only. Our method to learn explicit agent-agnostic navigational patterns is an alternative approach to enhance explainability of end-to-end learning, or incorporate an end-to-end learning aspect into the conventional modularized mobile robotics system [1].

III. PLAUSIBLE WORLD STATE INPUT GENERATION

Here we describe the pre-processing method shown in Fig. 1. Sensor observations are accumulated into partially observed world states x , which in turn are transformed into plausible world states \hat{x} . The proposed model uses \hat{x} as input.

A. Partial World State Representation

We generate partial world states based on accumulated sensor observations following the method described in prior work [15]. The method shares similarities with a hierarchical biological model of human representation and processing of visual information [41]. The agent is initialized within an unknown metric vector space. Sensor observations are projected onto this common vector space at discrete timesteps. Semantic information is inferred from images using a pretrained semantic segmentation model and appended to coincident 3D points to form semantic point clouds. Past semantic point clouds are integrated with new

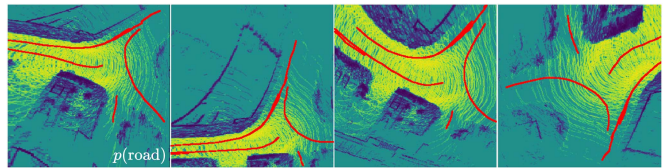


Fig. 2. Geometric data augmentation generates diverse sample variations from a single real sample. Spatial information (dense maps) and observed trajectories (red lines) are transformed by the same function.

observations by scan matching using the ICP algorithm [43] and SLAM [44] for loop closure. The accumulated semantic point cloud is reduced to a five-layered 2D probabilistic BEV representation $x \in \mathbb{R}^{I \times J \times C}$ with dimension $I \times J$ elements, and C denoting the number of semantic information channels. In this work, C consists of five channels representing the semantic attributes of a spatial point (i, j) ; we represent road probability $p(\text{road})$ by a beta distribution, lidar reflection intensity ϵ as a scalar value, and visual appearance by RGB values.

Dynamic objects are detected by a pretrained object detection model and represented by 3D bounding boxes. Trajectory observations are generated by temporally tracking detected objects. Dynamic objects are considered “moving” if motion is observed or “static” otherwise. This classification allows filtering away observations associated with moving dynamic objects while keeping observations of static dynamic objects for training, as they may influence how other agents navigate the environment such as swerving out of the lane to avoid a parked car. The static dynamic objects can be removed at inference time to provide an agent-agnostic prediction of navigational patterns akin to a lane map.

We leverage geometric data augmentation [16] to improve model generalization performance by learning geometric invariance. Each sample is augmented by random rotation and translation, and a polynomial warping function is applied to the dense maps and observed trajectories

$$a_0(\xi')^2 + a_1(\xi') + a_2 = \xi \quad (1)$$

where ξ is a substitute for spatial coordinates i and j , and ξ' denotes warped coordinates. We create dense warp maps by using the inverse function of (1) to map each warped coordinate ξ' to an original coordinate ξ . The coefficients a_0 , a_1 , and a_2 are derived by satisfying boundary conditions [16]. Fig. 2 shows visual examples of a sample augmentation.

B. Predictive World Model

The predictive world model [15] samples diverse and plausible complete world states \hat{x} conditioned on partially observed world states x as exemplified in Fig. 1. The world model is functionally similar to the biological ventral cortical pathway as the model disambiguates the partially observed environment by leveraging past experience [9]. The world model is computationally conceptualized as an arbitrary conditioning generative model and implemented by the recent SOTA hierarchical VAE (HVAE) model VDVAE [45] with the encoder module replaced by a posterior matching encoder [15]. The HVAE models the joint distribution of observable variables $p(r, \epsilon, R, G, B)$ factorized as the conditional distribution

$$p(r, \epsilon, R, G, B) = p(R|G, B, r)p(G|B, r)p(B|r)p(\epsilon|r)p(r) \quad (2)$$

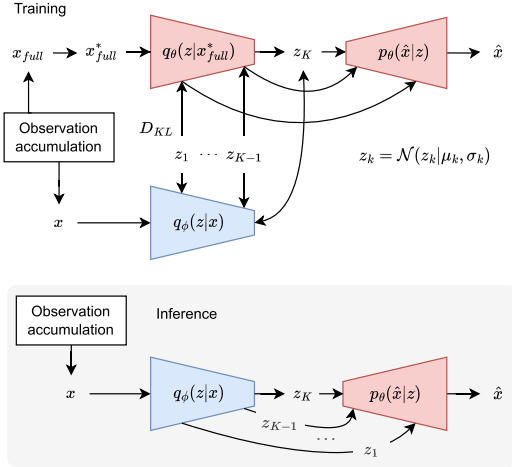


Fig. 3. (Top) Predictive world model is trained to reconstruct pseudo ground-truth world states x_{full}^* and simultaneously optimize a secondary encoder $q_\phi(z|x)$ to predict a similar hierarchical latent distribution z as the primary encoder $q_\theta(z|x_{full}^*)$ from partially observed world states x . (Bottom) The trained model samples a hierarchical latent distribution z from x to generate complete plausible world states \hat{x} .

using hierarchical latent variables z . Here r and ϵ denote road and lidar reflection intensity, and RGB are image color channels. The latent variable prior $p(z)$ and posterior $q(z|x)$ distributions are factorized as

$$p(z) = p(z_1|z_2) \dots p(z_{K-1}|z_K)p(z_K) \quad (3)$$

$$q(z|x) = q(z_1|z_2, x) \dots q(z_{K-1}|z_K, x)q(z_K|x) \quad (4)$$

with random variables z modeled by normal distributions.

The world model learns to approximate the prior and posterior distributions by the parameterized models $q_\theta(z|x)$ and $p_\theta(x|z)$ using variational inference [46] and trained using self-supervised learning to predict future observations from present observations akin to the predictive coding problem [47]. Note that the vanilla HVAE cannot learn to generate diverse complete representations from partially observed representations only. We follow the posterior matching optimization method visualized in Fig. 3 and presented in prior work [15] to overcome this limitation. The method trains a regular HVAE using pseudo ground-truth world states x_{full}^* , and a secondary encoder $q_\phi(z|x)$ to predict a similar hierarchical latent distribution $z = \{z_1, \dots, z_K\}$ as the primary encoder $q_\theta(z|x_{full}^*)$ from x .

We generate pseudo ground-truth world states x_{full}^* using a sequential process starting from the intermediate world state x_{full} consisting of past and future observations as explained in prior work [15]. The regular HVAE model is trained by maximizing the hierarchical ELBO over x_{full} [15], [45]. The second encoder is optimized by minimizing

$$\begin{aligned} & D_{KL}(q_\theta(z|x_{full}^*)||q_\phi(z|x)) \\ &= \sum_{k=1}^K \mathbb{E}_{q(z_{>k}|x)} [D_{KL}(q_\theta(z_k|z_{>k}, x_{full}^*)||q_\phi(z_k|z_{>k}, x))]. \end{aligned} \quad (5)$$

At inference time the model uses the partially observed encoder to generate a latent distribution $q_\phi(z|x)$ that can be decoded by $p_\theta(\hat{x}|z)$ into a completely observed plausible world state \hat{x}

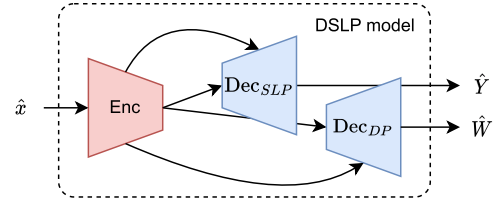


Fig. 4. Directional soft lane probability (DSL) model uses a dual decoder U-Net [48] model to transform a plausible world state \hat{x} into a soft lane probability (SLP) map \hat{Y} and directional probability (DP) tensor \hat{W} .

similar to a pseudo ground-truth world state x_{full}^* without the need to observe the future.

IV. DIRECTIONAL SOFT LANE PROBABILITY MODEL

Here we present a method to train a model to predict unbiased probability maps of local directional traversability. The model input is the plausible world state \hat{x} described in Section III. We also present a method for inferring global navigational patterns from the local probability maps. See Figs. 1 and 7 for output visualizations.

The model is implemented by a U-Net neural network [48] with a single encoder and two decoders as illustrated in Fig. 4. The first decoder outputs a probability map $Y \in \mathbb{R}^{I \times J}$ representing soft lane probabilities for elements in a grid map of size $I \times J$. The second decoder outputs a map of categorical distributions $W \in \mathbb{R}^{M \times I \times J}$ representing M direction interval probabilities for each location (i, j) . The methods for optimizing both probabilistic outputs are explained below.

A. Soft Lane Probability (SLP) Modeling

The likelihood of each environment location (i, j) being traversed by an unspecified agent is modeled by the predicted probability value $\hat{y}_{i,j} \in \hat{Y}$ and is called soft lane probability (SLP). Learning to predict an unbiased \hat{Y} from partial observations is nontrivial, as the self-supervised learning signal contains false negative traversal observations (i.e. lacking an observed trajectory where traversals are probable). We formalize the problem as follows. Ideally we want to learn a distribution $q(y)$ that approximates the true distribution $p(y)$. However, optimizing $q(y)$ according to the learning signal results in learning the distribution of partially observed samples $\tilde{p}(y)$. A principled solution is to use a regularizer to decrease bias and make $q(y)$ better match $p(y)$.

In this letter we present a semi-supervised objective that enables learning an unbiased probabilistic prediction of traversability based on an information-theoretic regularizer derived from balancing the information contribution from positive and negative partial observations in Y .

In information theory, the entropy $H(y)$ of a distribution $p(y)$ is considered a quantity that measures information content. The cross-entropy

$$H(p, q) \triangleq - \sum_{k=1}^K p(y = k) \log(q(y = k)) \quad (6)$$

measures the information overhead to compress a sample $y \sim p(y)$ using a code based on $q(y)$ [49].

Each partial observation Y contains two distinct groups of traversal information; a set of true positives representing certain information, and a set of true and false negatives representing uncertain information. The contributed information of the set of positive and negative observations are

$$H(Y_{pos} \subseteq Y, \hat{Y}) = - \sum_{i,j \in Y_{pos}} y_{i,j} \log(\hat{y}_{i,j}) \quad (7)$$

$$H(Y_{neg} \subseteq Y, \hat{Y}) = - \sum_{i,j \in Y_{neg}} (1 - y_{i,j}) \log(1 - \hat{y}_{i,j}). \quad (8)$$

We devise a regularizer based on balancing the information contribution provided by (7) and (8) according to the ratio of observations

$$\alpha_{IB} = |Y_{pos}| / (|Y_{pos}| + |Y_{neg}|) \quad (9)$$

where $|Y^*|$ denotes the number of positive and negative observed elements (i, j) . The balanced information contribution $H^*(Y|\hat{Y})$ is obtained by linearly interpolating the information contributions according to the ratio of observations

$$H^*(Y|\hat{Y}) = \alpha_{IB} H(Y_{neg}|\hat{Y}) + (1 - \alpha_{IB}) H(Y_{pos}|\hat{Y}). \quad (10)$$

Linear interpolation is a monotonic function that balances the information contributions while preserving the total information quantity

$$0 \leq H^*(Y|\hat{Y}) \leq \max(H(Y_{pos}|\hat{Y}), H(Y_{neg}|\hat{Y})). \quad (11)$$

We formulate the problem specific optimization objective \mathcal{L}_{SLP} as the mean balanced information contribution

$$\begin{aligned} \mathcal{L}_{SLP} = & - \frac{1}{|Y|} \sum_{i,j \in Y} [\alpha_{IB}(1 - y_{i,j}) \log(1 - \hat{y}_{i,j}) \\ & + (1 - \alpha_{IB}) y_{i,j} \log(\hat{y}_{i,j})] \end{aligned} \quad (12)$$

where $\hat{y}_{i,j}$ and $y_{i,j}$ is the predicted and observed soft lane probability for the element located at i, j . $|Y|$ denotes the number of traversable elements. The information contribution ratio α_{IB} provides the optimal interpolation between positive and negative traversal observations.

One can view (12) as the cross entropy objective with an additional dynamic regularizer between positive and negative observations. Experiments show that the balanced information contribution cross-entropy objective (12) performs better than finetuning a static hyperparameter weighting [16], and allows learning probabilistic predictions despite occasional abnormal observations unlike the barrier loss objective [31]. The negative log likelihood NLL_{SLP} of an observed sample y according to a model prediction \hat{y} based on modeling $p(y|\hat{y})$ as a Bernoulli distribution is

$$\text{NLL}_{SLP} = - \sum_{i,j \in Y} [y_{i,j} \log(\hat{y}_{i,j}) + (1 - y_{i,j}) \log(1 - \hat{y}_{i,j})]. \quad (13)$$

B. Directional Probability (DP) Modeling

The likelihood of local traversal directionality at each location (i, j) is modeled by the predicted vector $\hat{w}_{i,j}$ called directional probability (DP). The $\hat{w}_{i,j}$ models a categorical probability

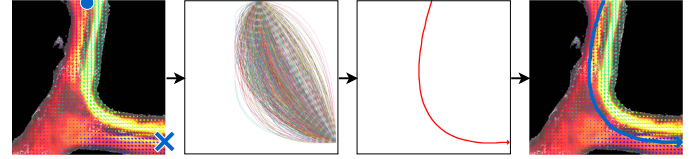


Fig. 5. Maximum likelihood graph is generated by connecting entry (\bullet) and exit (\times) points by the most probable of many sampled paths given the predicted DSLP field.

distribution representing the direction interval $\theta \in [0, 2\pi)$ by M uniformly spaced intervals

$$w_{i,j} = \left(p\left(\theta \in \left[0, \frac{2\pi}{M}\right)\right), \dots, p\left(\theta \in \left[\frac{(M-1)2\pi}{M}, 2\pi\right)\right) \right)^T. \quad (14)$$

The learning signal is created by encoding observed trajectories into $w_{i,j}$ as a discrete von Mises distribution. In the case of multiple overlapping trajectories the individual distributions are superimposed and renormalized. Learning to match distributions improve multimodal prediction compared with learning to predict single values by maximum likelihood estimation [16].

The optimization objective \mathcal{L}_{DP} is formulated as learning to predict the directional distribution by minimizing the mean KL divergence between predicted $\hat{w}_{i,j}$ and observed $w_{i,j}$ directionality over all elements $w_{i,j} \in W$

$$\mathcal{L}_{DP} = \frac{1}{|W|} \sum_{i,j \in W} D_{KL}(w_{i,j} || \hat{w}_{i,j}). \quad (15)$$

Note that the learning signal used to optimize the DP objective (15) lacks false negatives and therefore does not require regularization like the SLP objective (12).

The negative log likelihood NLL_{DP} of an observed sample $w_{i,j}$ according to a model prediction $\hat{w}_{i,j}$ based on modeling $p(w|\hat{w})$ as a categorical distribution is

$$\text{NLL}_{DP} = - \sum_{i,j \in Y} \sum_{m=1}^M w_{i,j}^{(m)} \log(\hat{w}_{i,j}^{(m)}). \quad (16)$$

C. Maximum Likelihood Lane Graph

Evaluating the goodness of local navigational patterns using the predicted DSLP field is straightforward. To also evaluate the usefulness of the predicted DSLP field for inferring global navigational patterns, we present a sampling-based method to generate a maximum likelihood road lane graph fitted to the predicted DSLP field. The graph generation process is illustrated in Fig. 5.

First, we infer entry and exit points at the edges of the predicted DSLP field. A non-maximum suppression (NMS) operation is performed on the SLP field \hat{Y} to find the most likely path centers. Each point is designated as an entry and/or exit point according to the predicted DP field \hat{W} . Additional entry and exit points are inferred from directional field regions which are coherent but lack a NMS point.

Secondly, we incrementally build a graph by searching for valid connecting paths between all entrance and exit points by a sampling-based approach. A set of second-degree polynomial spline paths is generated between an entry and exit pair by randomly sampling a valid spline control point $(i, j)^*$ from a

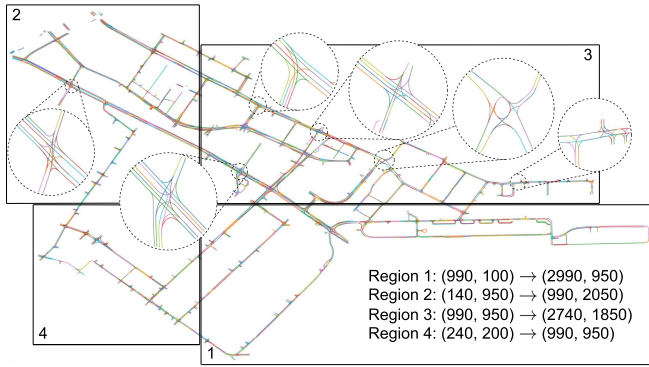


Fig. 6. Samples are partitioned into four nonoverlapping regions. Regions are specified by bottom-left and top-right corners in world coordinates.

normal distribution with rejection sampling. The likelihood of each sampled path is evaluated using the location and directionality of M equidistant points along the path given the predicted DSLP field using (13) and (16). The path with the lowest total NLL is selected as the best path. Repeating this process results in a set of most likely paths representing the maximum likelihood graph. A post-processing operation removes undesired edges between neighboring lanes (i.e. u-turns) using a simple distance threshold heuristic. Representing navigational patterns by splines is a useful inductive bias, as agents tend to navigate structured environments in a continuous and smooth manner.

V. EXPERIMENTS

We evaluate the model performance on the right-side driving daytime Boston scenes in the nuScenes dataset [50] similar to our baseline methods [6], [7]. The observation accumulation method described in Section III-A generates a partially observed training sample x every 1 m using accumulated observations from six 360° FoV RGB cameras and a top-mounted 32 beam lidar and a single pretrained semantic segmentation model [15]. Each x is augmented 20 times. Partitioning the generated training samples into the nonoverlapping regions shown in Fig. 6 results in 60,960 (34.7%), 40,960 (23.3%), and 73,780 (42.0%) samples for regions 1 to 3. Evaluation region 4 contains samples generated every 10 m without augmentation. We use a semantic segmentation model pretrained on two different public datasets [15]. We accumulate observations using ground truth pose information to reduce engineering effort, as prior work demonstrates the feasibility of accumulation based on pose estimation [15]. The plausible world state model input representation \hat{x} consists of a five-layered 256×256 grid map encompassing a 51.2×51.2 m region similar to prior work [6].

We conduct a model hyperparameter study and find that a smaller 1.4 M parameter model generalizes best. The model as depicted in Fig. 4 has a common 8-layered CNN encoder with filter count increasing from 16 to 256, and two 8-layered CNN decoders with bilinear upsampling and filter count decreasing from 64 to 8. See the code for further implementation details. We use the following benchmarks to evaluate our DSLP model. We compare the global navigation pattern inference performance against the two most relevant and recently published SOTA supervised models STSU [7] and LaneGraphNet [6]. Both baselines are trained on nuScenes data [50] to predict lane graphs using complete ground truth graphs as supervision. We compare

TABLE I
PERFORMANCE OF PREDICTED LOCAL PROBABILITY FIELDS

	NLL _{SLP}	NLL _{DP}	NLL	Dir. acc.	
DSLA [16]	2.499	12.596	15.095	0.864	
DSLSP	const α	0.423	12.241	12.663	0.855
	mean α_{IB}	0.444	12.038	12.482	0.881
	α_{IB}	0.556	11.769	12.325	0.892
	full obs.	0.539	11.666	12.205	0.900

the local probability field estimation performance against the prior self-supervised SOTA model called DSLA [16].

Local probability field estimation: We evaluate the predicted soft lane \hat{Y} and directional \hat{W} probability fields by computing the summed negative log-likelihood (NLL) of the ground truth lane map using (13) and (16). Lower NLL means the ground truth lane map is more likely according to the model. Directional accuracy measures the ratio of elements within $\pm 45^\circ$ of the ground truth direction.

Global navigational pattern inference: We evaluate the usefulness of the predicted probability fields for inferring global navigational patterns by computing the intersection over union (IoU) and F1 score between the maximum likelihood graph and ground truth lane map. Our method does not consider the spacing of graph nodes as an integral part of navigational patterns and thus does not view node displacement as a relevant performance metric.

Ablation studies: We evaluate the advantage of our proposed predictive world modeling approach [15] for learning navigational patterns from sampled plausible completed worlds \hat{x} instead of partially observed worlds x . We conduct an experiment using unaugmented samples to quantify the performance contribution of our geometric data augmentation method [16] on real-world data. We conduct experiments on dataset splits including a different number of regions to estimate how performance increases with additional data.

VI. RESULTS

Local probability field estimation: Table I presents evaluation results for the predicted probability fields. Our proposed DSLSP model optimized with the information balance regularizer α_{IB} (9) predicts the least biased probability field among all models trained and evaluated on accumulated past observation inputs. We conclude that the probabilistic objective (12) substantially reduces bias compared with the non-probabilistic DSLA affordance objective [16]. Training and evaluating on accumulated past and future observation inputs in an offline map creation manner (i.e. full obs.) reduces bias, demonstrating that more comprehensively observed environments result in better performance. We performed experiments with different constant α values to demonstrate the merit of the proposed hyperparameter-free regularizer α_{IB} (9). The best constant weight α value 0.1, found over five hyperparameter experiments, results in worse performance than using α_{IB} . We demonstrate the merit of dynamic, per-sample computed α_{IB} values (9) by running an experiment with the constant mean α_{IB} value 0.122 computed over all training samples, which results in worse performance. See Fig. 7 for probability field visualizations.

Global navigational pattern inference: Table II presents results showing that the maximum likelihood graph fitted to the probability field predicted by our self-supervised DSLSP

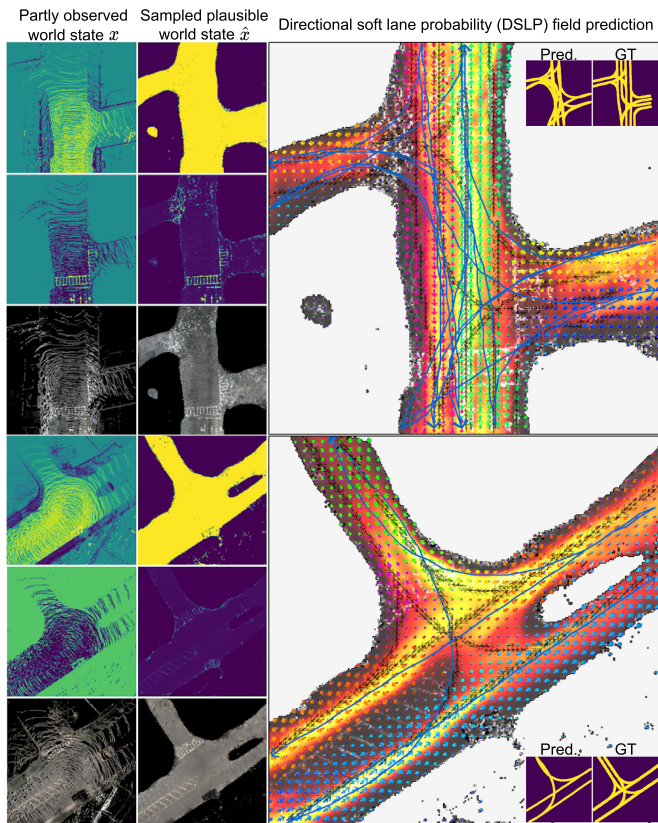


Fig. 7. Model output visualizations. The left column shows accumulated partial observations x . The middle column shows plausible world states \hat{x} sampled from x . The right column visualizes the predicted probability fields \hat{Y} and \hat{W} as well as the maximum likelihood graph.

TABLE II
PERFORMANCE OF GLOBAL NAVIGATIONAL PATTERN INFERENCE

	IoU	F1 score
STSU [7]	0.389	0.560
LaneGraphNet [6]	0.420	0.574
DSLPA [16]	0.427 (0.128)	0.839 (0.07)
DSLPL	constant α	0.418 (0.146) 0.853 (0.08)
	mean α_{IB}	0.410 (0.147) 0.846 (0.08)
	α_{IB}	0.442 (0.125) 0.834 (0.07)
	full obs.	0.454 (0.128) 0.839 (0.08)

and prior DSLPA model [16] from partially observed world representations x , outperforms the supervised SOTA baselines STSU [7] and LaneGraphNet [6] trained on ground truth lane graphs. Our self-supervised method not only improves upon the supervised baseline results while limited to the same training data domain, but is also a scalable solution for real-world mobile robotics as the model can improve by continual learning from new observational experience. While the baselines do not specify train and evaluation regions for an ideal comparison, our experiments in Table IV show our model surpassing the supervised baseline methods also when training on one region only, demonstrating that the exact train and evaluation region split is not critical for achieving our favorable results. We note that the probabilistic DSLPL model outperforms the non-probabilistic DSLPA affordance model [16], the proposed regularizer α_{IB} (9) outperforms the best constant hyperparameter regularizer α and the mean α_{IB} value, and that more comprehensively observed

TABLE III
ABLATION STUDIES

WM	Aug.	NLL _{SLP} *	NLL _{DP}	NLL*	Dir. acc.	IoU
✓	✓	0.189	11.769	11.958	0.892	0.442
✗	✓	0.266	12.785	13.051	0.853	0.223
✓	✗	0.167	13.764	13.931	0.848	0.453

*Mean over all elements

TABLE IV
PERFORMANCE WITH VARYING DATA AMOUNTS

# Regions	NLL _{SLP}	NLL _{DP}	NLL	Dir. acc.	IoU
{1}	0.478	12.696	13.174	0.861	0.423
{1, 2}	0.544	12.013	12.557	0.874	0.444
{1, 2, 3}	0.556	11.769	12.325	0.892	0.442

environments result in better performance. See Fig. 7 for inferred navigational path visualizations.

Ablation studies: Table III shows that leveraging the predictive world model (WM) [15] and proposed data augmentation (Aug.) [16] method reduces bias in the predicted probabilistic fields. We note that the unaugmented experiment generates output biased towards ego-agent trajectories, resulting in worse overall NLL while the maximum likelihood graph remains accurate. We believe this indicates the potential to further improve the graph generation algorithm to better leverage the more accurate probability field prediction. We do not explicitly evaluate the performance of the world model itself as this is done in prior work [15].

Table IV shows that increased observational experience reduces bias in the predicted probability field, providing evidence that the model can be trained to infer an unbiased probability prediction in the limit of infinite data

Inference time: We analyze the time taken for one iteration of our proposed system as follows. The mean inference time for the predictive world model and DSLPL model is 0.175 sec and 0.017 sec, resulting in a total mean time of 0.192 sec per iteration or 5.21 Hz on an RTX 4090 GPU. We conclude that our method is feasible to run in real-time as it introduces a 0.192 sec overhead with a real-time SLAM implementation [43] operating faster than sensor frame rates.

VII. CONCLUSION

In this letter, we present the first SSL method for training a model to infer navigational patterns in real-world environments from partial observations while achieving better performance than SOTA supervised baselines. Here we identify limitations and directions for future work. The representation of spatially small but semantically important environmental cues, such as road markings, is inefficiently represented by uniform grid maps. Traffic information on signs is not represented at all. We propose to instead detect and semantically draw road markings and signs in the input representation. Graph generation can be improved by inferring start and end points within the BEV, sampling higher-order splines, and decomposing splines into a sparse graph [31]. Understanding navigational patterns may require a temporal memory of past observations to resolve ambiguity. We propose an additional module that maintains a latent environment encoding by learning from sequences instead of i.i.d. data.

ACKNOWLEDGMENT

The authors would like to take this opportunity to thank the “Interdisciplinary Frontier Next-Generation Researcher Program of the Tokai Higher Education and Research System”.

REFERENCES

- [1] T. Krause, A. Pandey, R. Alami, and A. Kirsch, “Human-aware robot navigation: A survey,” *Robot. Auton. Syst.*, vol. 61, no. 12, pp. 1726–1743, 2013.
- [2] H. Sheif and X. Hu, “Autonomous driving in the iCity - HD maps as a key challenge of the automotive industry,” *Engineering*, vol. 2, pp. 159–162, 2016.
- [3] C. J. Holder and T. P. Breckong, “Learning to drive: End-to-end off-road path prediction,” *IEEE Intell. Transp. Syst. Mag.*, vol. 13, no. 2, pp. 217–221, Summer 2021, doi: [10.1109/MITS.2019.2898970](https://doi.org/10.1109/MITS.2019.2898970).
- [4] A. Amini, G. Rosman, S. Karaman, and D. Rus, “Variational end-to-end navigation and localization,” in *Proc. IEEE Int. Conf. Robot. Automat.*, 2018, pp. 8958–8964.
- [5] M. Bansal, A. Krizhevsky, and A. Ogale, “ChauffeurNet: Learning to drive by imitating the best and synthesizing the worst,” in *Proc. Robot.: Sci. Syst.*, 2019. [Online]. Available: <https://www.roboticsproceedings.org/rss15/p31.pdf>
- [6] J. Zürn, J. Vertens, and W. Burgard, “Lane graph estimation for scene understanding in urban driving,” *IEEE Robot. Automat. Lett.*, vol. 6, no. 4, pp. 8615–8622, Oct. 2021.
- [7] Y. B. Can, A. Liniger, D. P. Paudel, and L. Van Gool, “Structured bird’s-eye-view traffic scene understanding from onboard images,” in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2021, pp. 15661–15670.
- [8] J. Gibson, *The Ecological Approach to Visual Perception*. Boston, MA, USA: Houghton Mifflin, 1979.
- [9] D. Milner and M. Goodale, “Two visual systems re-viewed,” *Neuropsychologia*, vol. 46, no. 3, pp. 774–785, 2008.
- [10] Z. Han and A. Sereno, “Modeling the ventral and dorsal cortical visual pathways using artificial neural networks,” *Neural Computation*, vol. 34, no. 1, pp. 138–171, 2022.
- [11] D. Milner, “Is visual processing in the dorsal stream accessible to consciousness?,” *Proc. Roy. Soc. B: Biol. Sci.*, vol. 279, pp. 2289–2298, 2012.
- [12] B. Sheth and R. Young, “Two visual pathways in primates based on sampling of space: Exploitation and exploration of visual information,” *Front. Integr. Neurosci.*, vol. 10, 2016, Art. no. 37.
- [13] T. Lesort, V. Lomonaco, A. Stoian, D. Maltoni, and D. Filliat, “Continual learning for robotics: Definition, framework, learning strategies, opportunities and challenges,” *Inf. Fusion*, vol. 58, pp. 52–68, 2020.
- [14] V. Mnih et al., “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, pp. 529–533, 2015.
- [15] R. Karlsson, A. Carballo, K. Fujii, K. Ohtani, and K. Takeda, “Predictive world models from real-world partial observations,” in *Proc. IEEE Int. Conf. Mobility Operations Serv. Technol.*, 2023.
- [16] R. Karlsson and E. Sjöberg, “Learning a directional soft lane affordance model for road scenes using self-supervision,” in *Proc. IEEE Intell. Veh. Symp.*, 2020, pp. 1141–1148.
- [17] T. Salzmann, J. Thomas, T. Kühbeck, J. c. Sung, S. Wagner, and A. Knoll, “Online path generation from sensor data for highly automated driving functions,” in *Proc. IEEE Intell. Transp. Syst. Conf.*, 2019, pp. 1807–1812.
- [18] U. Baumann, C. Guiser, M. Herman, and J. Zollner, “Predicting ego-vehicle paths from environmental observations with a deep neural network,” in *Proc. IEEE Int. Conf. Robot. Automat.*, 2018, pp. 4709–4716.
- [19] D. Barnes, W. Maddern, and I. Posner, “Find your own way: Weakly supervised segmentation of path proposals for urban autonomy,” in *Proc. IEEE Int. Conf. Robot. Automat.*, 2017, pp. 203–210.
- [20] T. Ort et al., “MapLite: Autonomous intersection navigation without a detailed prior map,” *IEEE Robot. Automat. Lett.*, vol. 5, no. 2, pp. 556–563, Apr. 2020.
- [21] S. Casas, A. Sadat, and R. Urtasun, “MP3: A unified model to map, perceive, predict, and plan,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 14398–14407.
- [22] N. Pérez-Higueras, F. Caballero, and L. Merino, “Learning human-aware path planning with fully convolutional networks,” in *Proc. IEEE Int. Conf. Robot. Automat.*, 2018, pp. 5897–5902.
- [23] S. Karaman and E. Frazzoli, “Sampling-based algorithms for optimal motion planning,” *Int. J. Robot. Res.*, vol. 30, no. 7, pp. 846–894, 2011.
- [24] K. Kitani, B. Ziebart, A. Bagnell, and M. Hebert, “Activity forecasting,” in *Proc. Eur. Conf. Comput. Vis.*, 2012, pp. 201–214.
- [25] N. Ratliff, D. Silver, and J. Bagnell, “Learning to search: Functional gradient techniques for imitation learning,” *Auton. Robots*, vol. 27, no. 1, pp. 25–53, 2009.
- [26] N. Homayounfar, W. C. Ma, S. K. Lakshmikanth, and R. Urtasun, “Hierarchical recurrent attention networks for structured online maps,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 3417–3426.
- [27] N. Homayounfar et al., “DAGMapper: Learning to map by discovering lane topology,” in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2019, pp. 2911–2920.
- [28] Y. Guo et al., “Gen-LaneNet: A generalized and scalable approach for 3D lane detection,” in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 666–681.
- [29] L. Zhang et al., “Hierarchical road topology learning for urban map-less driving,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2022, pp. 3816–3823.
- [30] L. Mi et al., “HDMaGen: A hierarchical graph generative model of high definition maps,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 4225–4234.
- [31] R. Karlsson, D. Wong, S. Thompson, and K. Takeda, “Learning a model for inferring a spatial road lane network graph using self-supervision,” in *Proc. IEEE Int. Intell. Transp. Syst. Conf.*, 2021, pp. 812–819.
- [32] Q. Li, Y. Wang, Y. Wang, and H. Zhao, “HDMaNet: An online HD map construction and evaluation framework,” in *Proc. IEEE Int. Conf. Robot. Automat.*, 2022, pp. 4628–4634.
- [33] Y. Liu, Y. Yuan, Y. Wang, Y. Wang, and H. Zhao, “VectorMapNet: End-to-end vectorized HD map learning,” in *Proc. Int. Conf. Mach. Learn.*, 2023, pp. 22352–22369.
- [34] B. Liao et al., “MapTR: Structured modeling and learning for online vectorized HD map construction,” in *Proc. 11th Int. Conf. Learn. Representations*, 2023. [Online]. Available: https://openreview.net/forum?id=k7p_YAO7yE
- [35] J. Shin, F. Rameau, H. Jeong, and D. Kum, “InstaGraM: Instance-level graph modeling for vectorized HD map learning,” in *Proc. Vis.-Centric Auton. Driving Workshop*, 2023. [Online]. Available: <https://vcad.site/#/>
- [36] T. Ort, J. Walls, S. A. Parkinson, I. Gilitschenski, and D. Rus, “MapLite 2.0: Online HD map inference using a prior SD map,” *IEEE Robot. Automat. Lett.*, vol. 7, no. 3, pp. 8355–8362, Jul. 2022.
- [37] D. Pomerleau, “ALVINN: An autonomous land vehicle in a neural network,” in *Proc. Neural Inf. Process. Syst.*, 1988, pp. 305–313.
- [38] A. Kendall et al., “Learning to drive in a day,” in *Proc. IEEE Int. Conf. Robot. Automat.*, 2019, pp. 8248–8254.
- [39] M. Henaff, A. Canziani, and Y. LeCun, “Model-predictive policy learning with uncertainty regularization for driving in dense traffic,” in *Proc. Int. Conf. Learn. Representations*, 2019. [Online]. Available: <https://openreview.net/forum?id=HygQBn0cYm>
- [40] D. Chen, V. Koltun, and P. Krähenbühl, “Learning to drive from a world on rails,” in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2021, pp. 15590–15599.
- [41] D. Marr, *Vision: A Computational Investigation Into the Human Representation and Processing of Visual Information*, W. H. Freeman. Cambridge, MA, USA: MIT Press, 1982.
- [42] P. J. Besl and N. D. McKay, “A method for registration of 3-D shapes,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 14, no. 2, pp. 239–256, Feb. 1992.
- [43] I. Vizzo, T. Guadagnino, B. Mersch, L. Wiesmann, J. Behley, and C. Stachniss, “KISS-ICP: In defense of point-to-point ICP,” *IEEE Robot. Automat. Lett.*, vol. 8, no. 2, pp. 1029–1036, Feb. 2023.
- [44] R. Smith and P. Cheeseman, “On the representation and estimation of spatial uncertainty,” *Int. J. Robot. Res.*, vol. 5, no. 4, pp. 56–68, 1986.
- [45] R. Child, “Very deep VAEs generalize autoregressive models and can outperform them on images,” in *Proc. Int. Conf. Learn. Representations*, 2021. [Online]. Available: <https://openreview.net/forum?id=RLRXCV6DbEJ>
- [46] D. Kingma and M. Welling, “Auto-encoding variational Bayes,” in *Proc. Int. Conf. Learn. Representations*, 2014. [Online]. Available: <https://openreview.net/forum?id=33X9fd2-9FyZd>
- [47] J. Marino, “Predictive coding, variational autoencoders, and biological connections,” *Neural Computation*, vol. 34, pp. 1–44, 2019.
- [48] O. Ronneberger, P. Fischer, and T. Brox, “U-Net: Convolutional networks for biomedical image segmentation,” in *Proc. 18th Med. Image Comput. Comput.-Assist. Interv.: 18th Int. Conf.*, 2015, pp. 234–241.
- [49] K. Murphy, *Probabilistic Machine Learning*. Cambridge, MA, USA: MIT Press, 2022.
- [50] H. Caesar et al., “nuScenes: A multimodal dataset for autonomous driving,” in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 11621–11631.